# Process/Thread Pinning

Pinning, the binding of a process or thread to a specific core, can improve the performance of your code.

You can insert the MPI function mpi_get_processor_name and the Linux C function sched_getcpu into your source code to check process and/or thread placement. The MPI function mpi_get_processor_name returns the hostname an MPI process is running on (to be used for MPI and/or MPI+OpenMP codes only). The Linux C function sched_getcpu returns the processor number the process/thread is running on.

If your source code is written in Fortran, you can use the C codemycpu.c, which allows your Fortran code to call sched_getcpu. The next section describes how to use the mycpu.c code.

## C Program mycpu.c

```
#include <utmpx.h>
int sched_getcpu();

int findmycpu_ ()
{
    int cpu;
    cpu = sched_getcpu();
    return cpu;
}
```

Compile mycpu.c as follows to produce the object file mycpu.o:

```
pfe21% module load comp-intel/2020.4.304
pfe21% icc -c mycpu.c
```

The following example demonstrates how to instrument an MPI+OpenMP source code with the above functions. The added lines are highlighted.

```
    program your_program
    use omp_lib
...
    integer :: resultlen, tn, cpu
    integer, external :: findmycpu
    character (len=8) :: name

    call mpi_init( ierr )
    call mpi_comm_rank( mpi_comm_world, rank, ierr )
    call mpi_comm_size( mpi_comm_world, numprocs, ierr )
    call mpi_get_processor_name(name, resultlen, ierr)
!$omp parallel

    tn = omp_get_thread_num()
    cpu = findmycpu()
    write (6,*) 'rank ', rank, ' thread ', tn,
    &  ' hostname ', name, ' cpu ', cpu
.....
!$omp end parallel
    call mpi_finalize(ierr)
    end
```

Compile your instrumented code as follows:

```
pfe21% module load comp-intel/2020.4.304
pfe21% module load mpi-hpe/mpt
pfe21% ifort -o a.out -qopenmp mycpu.o your_program.f -lmpi
```

## Sample PBS script

The following PBS script provides an example of running the hybrid MPI+OPenMP code across two nodes, with 2 MPI processes per node and 4 OpenMP threads per process, and using the [mbind](#) tool to pin the processes and threads.

```
#PBS -lselect=2:ncpus=28:mpiprocs=2:model=bro
#PBS -lwalltime=0:10:00

cd $PBS_O_WORKDIR

module load comp-intel/2020.4.304
module load mpi-hpe/mpt

mpiexec -np 4 mbind.x -cs -t4 -v ./a.out
```

Here is a sample output:

These 4 lines are generated by mbind only if you have included the -v option:

```
host: r627i4n1, ncpus: 56, rank: 0 (r0), nthreads: 4, bound to cpus: {0-9:3}
host: r627i4n1, ncpus: 56, rank: 1 (r1), nthreads: 4, bound to cpus: {14-23:3}
host: r627i4n8, ncpus: 56, rank: 2 (r0), nthreads: 4, bound to cpus: {0-9:3}
host: r627i4n8, ncpus: 56, rank: 3 (r1), nthreads: 4, bound to cpus: {14-23:3}
```

These lines are generated by your instrumented code:
```
rank   0 thread   3 hostname r627i4n1 cpu   0
rank   0 thread   3 hostname r627i4n1 cpu   3
rank   0 thread   3 hostname r627i4n1 cpu   6
rank   0 thread   3 hostname r627i4n1 cpu   9
rank   1 thread   3 hostname r627i4n1 cpu   14
rank   1 thread   3 hostname r627i4n1 cpu   17
rank   1 thread   3 hostname r627i4n1 cpu   20
rank   1 thread   3 hostname r627i4n1 cpu   23
rank   2 thread   3 hostname r627i4n8 cpu   0
rank   2 thread   3 hostname r627i4n8 cpu   3
rank   2 thread   3 hostname r627i4n8 cpu   6
rank   2 thread   3 hostname r627i4n8 cpu   9
rank   3 thread   3 hostname r627i4n8 cpu   14
rank   3 thread   3 hostname r627i4n8 cpu   17
rank   3 thread   3 hostname r627i4n8 cpu   20
rank   3 thread   3 hostname r627i4n8 cpu   23
```

Note: In your output, these lines may be listed in a different order.

# Using Intel OpenMP Thread Affinity for Pinning

UPDATE IN PROGRESS: This article is being updated to support Skylake and Cascade Lake.

The Intel compiler's OpenMP runtime library has the ability to bind OpenMP threads to physical processing units. Depending on the system topology, application, and operating system, thread affinity can have a dramatic effect on code performance. We recommend two approaches for using the Intel OpenMP thread affinity capability.

## Using the KMP_AFFINITY Environment Variable

The thread affinity interface is controlled using the KMP_AFFINITY environment variable.

## Syntax

For csh and tcsh:

setenv KMP_AFFINITY [<*modifier*>,...]<*type*>[,<*permute*>][,<*offset*>]

For sh, bash, and ksh:

export KMP_AFFINITY=[<*modifier*>,...]<*type*>[,<*permute*>][,<*offset*>]

## Using the Compiler Flag -par-affinity Compiler Option

Starting with the Intel compiler version 11.1, thread affinity can be specified through the compiler option -par-affinity. The use of -openmp or -parallel is required in order for this option to take effect. This option overrides the environment variable when both are specified. See **man ifort** for more information.

Note: Starting with comp-intel/2015.0.090, -openmp is deprecated and has been replaced with -qopenmp.

## Syntax

-par-affinity=[<*modifier*>,...]<*type*>[,<*permute*>][,<*offset*>]

## Possible Values for type

For both of the recommended approaches, the only required argument is type, which indicates the type of thread affinity to use. Descriptions of all of the possible arguments (type, modifier, permute, and offset) can be found in man ifort.

**Recommendation:** Use Intel compiler versions 11.1 and later, as some of the affinity methods described below are not supported in earlier versions.

Possible values for type are:

type = none (default)
   Does not bind OpenMP threads to particular thread contexts; however, if the operating system supports affinity, the compiler still uses the OpenMP thread affinity interface to determine machine topology. Specify KMP_AFFINITY=verbose,none to list a machine topology map.
type = disabled
   Specifying disabled completely disables the thread affinity interfaces. This forces the OpenMP runtime library to behave as if the affinity interface was not supported by the operating system. This includes implementations of the low-level API interfaces such as kmp_set_affinity and kmp_get_affinity that have no effect and will return a nonzero error code.

Additional information from Intel:

"The thread affinity type of KMP_AFFINITY environment variable defaults to none (KMP_AFFINITY=none). The behavior for KMP_AFFINITY=none was changed in 10.1.015 or later, and in all 11.$x$ compilers, such that the initialization thread creates a "full mask" of all the threads on the machine, and every thread binds to this mask at startup time. It was subsequently found that this change may interfere with other platform affinity mechanism, for example, dplace() on Altix machines. To resolve this issue, a new affinity type disabled was introduced in compiler 10.1.018, and in all 11.$x$ compilers (KMP_AFFINITY=disabled). Setting KMP_AFFINITY=disabled will prevent the runtime library from making any affinity-related system calls."

type = compact
   Specifying compact causes the threads to be placed as close together as possible. For example, in a topology map, the nearer a core is to the root, the more significance the core has when sorting the threads.

Usage example:

```
# for sh, ksh, bash
export KMP_AFFINITY=compact,verbose

# for csh, tcsh
setenv KMP_AFFINITY compact,verbose
```

type = scatter
>    Specifying scatter distributes the threads as evenly as possible across the entire system. Scatter is the opposite of compact.

Note: For most OpenMP codes, type=scatter should provide the best performance, as it minimizes cache and memory bandwidth contention for all processor models.

Usage example:

```
# for sh, ksh, bash
export KMP_AFFINITY=scatter,verbose

# for csh, tcsh
setenv KMP_AFFINITY scatter,verbose
```

type = explicit
>    Specifying explicit assigns OpenMP threads to a list of OS proc IDs that have been explicitly specified by using the proclist=*modifier*, which is required for this affinity type.

Usage example:

```
# for sh, ksh, bash
export KMP_AFFINITY="explicit,proclist=[0,1,4,5],verbose"

# for csh, tcsh
setenv KMP_AFFINITY "explicit,proclist=[0,1,4,5],verbose"
```

For nodes that support hyperthreading, you can use the granularity modifier to specify whether to pin OpenMP threads to physical cores using granularity=core (the default) or pin to logical cores using granularity=fine or granularity=thread for the compact and scatter types.

## Examples

The following examples illustrate the thread placement of an OpenMP job with four threads on various platforms with different thread affinity methods. The variable OMP_NUM_THREADS is set to 4:

```
# for sh, ksh, bash
export OMP_NUM_THREADS=4

# for csh, tcsh
setenv OMP_NUM_THREADS 4
```

The use of the verbose modifier is recommended, as it provides an output with the placement.

## Sandy Bridge (Pleiades)

As seen in the configuration diagram of a Sandy Bridge node, each set of eight physical cores in a socket share the same L3 cache.
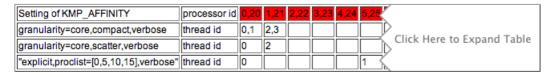
Four threads running on 1 node (16 physical cores and 32 logical cores due to hyperthreading) of Sandy Bridge will get the following thread placement:

| setting of KMP_AFFINITY | Processor id | 0,16 | 1,17 | 2,18 | 3,19 | 4,20 | 5,21 | |
|---|---|---|---|---|---|---|---|---|
| granularity=core,compact,verbose | thread id | 0,1 | 2,3 | | | | | |
| granularity=core,scatter,verbose | thread id | 0 | 2 | | | | | Click Here to Expand Table |
| "explicit,proclist=[0,4,8,12],verbose" | thread id | 0 | | | | 1 | | |

## Ivy Bridge (Pleiades)

As seen in the configuration diagram of an Ivy Bridge node, each set of ten physical cores in a socket share the same L3 cache.

Four threads running on 1 node (20 physical cores and 40 logical cores due to hyperthreading) of Ivy Bridge will get the following thread placement:

| Setting of KMP_AFFINITY | processor id | 0,20 | 1,21 | 2,22 | 3,23 | 4,24 | 5,25 | |
|---|---|---|---|---|---|---|---|---|
| granularity=core,compact,verbose | thread id | 0,1 | 2,3 | | | | | |
| granularity=core,scatter,verbose | thread id | 0 | 2 | | | | | Click Here to Expand Table |
| "explicit,proclist=[0,5,10,15],verbose" | thread id | 0 | | | | | 1 | |

## Haswell (Pleiades)

As seen in the configuration diagram of a [Haswell](#) node, each set of 12 physical cores in a socket share the same L3 cache.

Four threads running on 1 node (24 physical cores and 48 logical cores due to hyperthreading) of Haswell will get the following thread placement:

| Setting of KMP_AFFINITY | | 0,24 | 1,25 | 2,26 | 3,27 | 4,28 | 5,29 | |
|---|---|---|---|---|---|---|---|---|
| granularity=core,compact,verbose | thread id | 0,1 | 2,3 | | | | | Click Here to Expand Table |
| granularity=core,scatter,verbose | thread id | 0 | 2 | | | | | |
| "explicit,proclist=[0,6,12,18],verbose" | thread id | 0 | | | | | | |

Note: "processor id" appears in the header column for the red row.

## Broadwell (Pleiades and Electra)

As seen in the configuration diagram of a [Broadwell](#) node, each set of 14 physical cores in a socket share the same L3 cache.

Four threads running on 1 node (28 physical cores and 56 logical cores due to hyperthreading) of Broadwell will get the following thread placement:

| Setting of KMP_AFFINITY | processor id | 0,28 | 1,29 | 2,30 | 3,31 | 4,32 | 5,33 | |
|---|---|---|---|---|---|---|---|---|
| granularity=core,compact,verbose | thread id | 0,1 | 2,3 | | | | | Click Here to Expand Table |
| granularity=core,scatter,verbose | thread id | 0 | 2 | | | | | |
| "explicit_proclist=[0,7,14,21],verbose" | thread id | 0 | | | | | | |

For MPI codes built with HPE's MPT libraries, one way to control pinning is to set certain MPT memory placement environment variables. For an introduction to pinning at NAS, see [Process/Thread Pinning Overview](#).

## MPT Environment Variables

Here are the MPT memory placement environment variables:

### MPI_DSM_VERBOSE

Directs MPI to display a synopsis of the NUMA and host placement options being used at run time to the standard error file.

Default: Not enabled

The setting of this environment variable is ignored if MPI_DSM_OFF is also set.

### MPI_DSM_DISTRIBUTE

Activates NUMA job placement mode. This mode ensures that each MPI process gets a unique CPU and physical memory on the node with which that CPU is associated. Currently, the CPUs are chosen by simply starting at relative CPU 0 and incrementing until all MPI processes have been forked.

Default: Enabled

WARNING: If the nodes used by your job are not fully populated with MPI processes, use MPI_DSM_CPULIST, dplace, or omplace for pinning instead of MPI_DSM_DISTRIBUTE.

The MPI_DSM_DISTRIBUTE setting is ignored if MPI_DSM_CPULIST is also set, or if dplace or omplace are used.

### MPI_DSM_CPULIST

Specifies a list of CPUs on which to run an MPI application, excluding the shepherd process(es) and mpirun. The number of CPUs specified should equal the number of MPI processes (excluding the shepherd process) that will be used.

Syntax and examples for the list:

- Use a comma and/or hyphen to provide a delineated list:

  ```
  # place MPI processes ranks 0-2 on CPUs 2-4
  # and ranks 3-5 on CPUs 6-8
  setenv MPI_DSM_CPULIST "2-4,6-8"
  ```

- Use a "/" and a stride length to specify CPU striding:

  ```
  # Place the MPI ranks 0 through 3 stridden
  # on CPUs 8, 10, 12, and 14
  setenv MPI_DSM_CPULIST 8-15/2
  ```

- Use a colon to separate CPU lists of multiple hosts:

  ```
  # Place the MPI processes 0 through 7 on the first host
  # on CPUs 8 through 15. Place MPI processes 8 through 15
  # on CPUs 16 to 23 on the second host.
  setenv MPI_DSM_CPULIST 8-15:16-23
  ```

- Use a colon followed by allhosts to indicate that the prior list pattern applies to all subsequent hosts/executables:

  ```
  # Place the MPI processes onto CPUs 0, 2, 4, 6 on all hosts
  setenv MPI_DSM_CPULIST 0-7/2:allhosts
  ```

## Examples

An MPI job requesting 2 nodes on Pleiades and running 4 MPI processes per node will get the following placements, depending on the environment variables set:

```
#PBS -lselect=2:ncpus=8:mpiprocs=4
module load <mpt_module>
setenv ....
cd $PBS_O_WORKDIR
mpiexec -np 8 ./a.out
```

- setenv MPI_DSM_VERBOSE
  setenv MPI_DSM_DISTRIBUTE

  ```
  MPI: DSM information
  MPI: MPI_DSM_DISTRIBUTE enabled
  grank  lrank  pinning  node name    cpuid
     0     0    yes      r86i3n5        0
     1     1    yes      r86i3n5        1
     2     2    yes      r86i3n5        2
     3     3    yes      r86i3n5        3
     4     0    yes      r86i3n6        0
     5     1    yes      r86i3n6        1
     6     2    yes      r86i3n6        2
     7     3    yes      r86i3n6        3
  ```

- setenv MPI_DSM_VERBOSE
  setenv MPI_DSM_CPULIST 0,2,4,6

  ```
  MPI: WARNING MPI_DSM_CPULIST CPU placement spec list is too short.
  MPI:      MPI processes on host #1 and later will not be pinned.
  MPI: DSM information
  grank  lrank  pinning  node name    cpuid
     0     0    yes      r22i1n7        0
     1     1    yes      r22i1n7        2
     2     2    yes      r22i1n7        4
     3     3    yes      r22i1n7        6
     4     0    no       r22i1n8        0
     5     1    no       r22i1n8        0
     6     2    no       r22i1n8        0
     7     3    no       r22i1n8        0
  ```

- setenv MPI_DSM_VERBOSE
  setenv MPI_DSM_CPULIST 0,2,4,6:0,2,4,6

  ```
  MPI: DSM information
  grank  lrank  pinning  node name    cpuid
     0     0    yes      r13i2n12       0
     1     1    yes      r13i2n12       2
     2     2    yes      r13i2n12       4
     3     3    yes      r13i2n12       6
     4     0    yes      r13i3n7        0
     5     1    yes      r13i3n7        2
     6     2    yes      r13i3n7        4
     7     3    yes      r13i3n7        6
  ```

- setenv MPI_DSM_VERBOSE
  setenv MPI_DSM_CPULIST 0,2,4,6:allhosts

  ```
  MPI: DSM information
  grank  lrank  pinning  node name    cpuid
     0     0    yes      r13i2n12       0
     1     1    yes      r13i2n12       2
     2     2    yes      r13i2n12       4
     3     3    yes      r13i2n12       6
     4     0    yes      r13i3n7        0
     5     1    yes      r13i3n7        2
     6     2    yes      r13i3n7        4
     7     3    yes      r13i3n7        6
  ```

HPE's omplace is a wrapper script for dplace. It pins processes and threads for better performance and provides an easier syntax than dplace for pinning processes and threads.

The omplace wrapper works with HPE MPT as well as with Intel MPI. In addition to pinning pure MPI or pure OpenMP applications, omplace can also be used for pinning hybrid MPI/OpenMP applications.

A few issues with omplace to keep in mind:

- dplace and omplace do not work with Intel compiler versions 10.1.015 and 10.1.017. Use the Intel compiler version 11.1 or later, instead
- To avoid interference between dplace/omplace and Intel's thread affinity interface, set the environment variable KMP_AFFINITY to disabled or set OMPLACE_AFFINITY_COMPAT to ON
- The omplace script is part of HPE's MPT, and is located under the /nasa/hpe/mpt/*mpt_version_number*/bin directory

## Syntax

**For OpenMP:**
setenv OMP_NUM_THREADS *nthreads*
omplace [OPTIONS] *program args...*
or
omplace -nt *nthreads* [OPTIONS] *program args...*

**For MPI:**
mpiexec -np *nranks* omplace [OPTIONS] *program args...*

**For MPI/OpenMP hybrid:**
setenv OMP_NUM_THREADS *nthreads*
mpiexec -np *nranks* omplace [OPTIONS] *program args...*
or
mpiexec -np *nranks* omplace -nt *nthreads* [OPTIONS] *program args...*

Some useful omplace options are listed below:

WARNING: For omplace, a blank space is required between -c and cpulist. Without the space, the job will fail. This is different from dplace.

-b *basecpu*
> Specifies the starting CPU number for the effective CPU list.

-c *cpulist*
> Specifies the effective CPU list. This is a comma-separated list of CPUs or CPU ranges.

-nt *nthreads*
> Specifies the number of threads per MPI process. If this option is unspecified, it defaults to the value set for the OMP_NUM_THREADS environment variable. If OMP_NUM_THREADS is not set, then *nthreads* defaults to 1.

-v
> Verbose option. Portions of the automatically generated placement file will be displayed.

-vv
> Very verbose option. The automatically generated placement file will be displayed in its entirety.

For information about additional options, see **man omplace**.

## Examples

## For Pure OpenMP Codes Using the Intel OpenMP Library

Sample PBS script:

#PBS -lselect=1:ncpus=12:model=wes

module load comp-intel/2015.0.090
setenv KMP_AFFINITY disabled

omplace -c 0,3,6,9 -vv ./a.out

Sample placement information for this script is given in the application's stout file:

omplace: placement file /tmp/omplace.file.21891
   firsttask cpu=0
   thread oncpu=0 cpu=3-9:3 noplace=1  exact

The above placement output may not be easy to understand. A better way to check the placement is to run the ps command on the running host while the job is still running:

ps -C a.out -L -opsr,comm,time,pid,ppid,lwp > placement.out

Sample output of placement.out

```
PSR COMMAND          TIME  PID  PPID  LWP
  0 openmp1        00:00:02 31918 31855 31918
 19 openmp1        00:00:00 31918 31855 31919
  3 openmp1        00:00:02 31918 31855 31920
  6 openmp1        00:00:02 31918 31855 31921
  9 openmp1        00:00:02 31918 31855 31922
```

Note that Intel OpenMP jobs use an extra thread that is unknown to the user, and does not need to be placed. In the above example, this extra thread is running on logical core number 19.

## For Pure MPI Codes Using HPE MPT

Sample PBS script:

#PBS -l select=2:ncpus=12:mpiprocs=4:model=wes

module load comp-intel/2015.0.090
module load mpi-hpe/mpt

#Setting MPI_DSM_VERBOSE allows the placement information
#to be printed to the PBS stderr file

setenv MPI_DSM_VERBOSE

mpiexec -np 8 omplace -c 0,3,6,9 ./a.out

Sample placement information for this script is shown in the PBS stderr file:

```
MPI: DSM information
MPI: using dplace
grank  lrank  pinning  node name     cpuid
  0      0    yes    r144i3n12        0
  1      1    yes    r144i3n12        3
  2      2    yes    r144i3n12        6
  3      3    yes    r144i3n12        9
  4      0    yes    r145i2n3         0
  5      1    yes    r145i2n3         3
  6      2    yes    r145i2n3         6
  7      3    yes    r145i2n3         9
```

In this example, the four processes on each node are evenly distributed to the two sockets (CPUs 0 and 3 are on the first socket while CPUs 6 and 9 on the second socket) to minimize contention. If omplace had not been used, then placement would follow the rules of the environment variable OMP_DSM_DISTRIBUTE, and all four processes would have been placed on the first socket -- likely leading to more contention.

## For MPI/OpenMP Hybrid Codes Using HPE MPT and Intel OpenMP

Proper placement is more critical for MPI/OpenMP hybrid codes than for pure MPI or pure OpenMP codes. The following example demonstrates the situation when no placement instruction is provided and the OpenMP threads for each MPI process are stepping on one another which likely would lead to very bad performance.

Sample PBS script without pinning:

#PBS -l select=2:ncpus=12:mpiprocs=4:model=wes

module load comp-intel/2015.0.090
module load mpi-hpe/mpt
setenv OMP_NUM_THREADS 2

mpiexec -np 8 ./a.out

There are two problems with the resulting placement shown in the example above. First, you can see that the first four MPI processes on each node are placed on four cores (0,1,2,3) of the same socket, which will likely lead to more contention compared to when they are distributed between the two sockets.

```
MPI: MPI_DSM_DISTRIBUTE enabled
grank  lrank  pinning  node name     cpuid
  0      0    yes    r212i0n10        0
  1      1    yes    r212i0n10        1
  2      2    yes    r212i0n10        2
  3      3    yes    r212i0n10        3
  4      0    yes    r212i0n11        0
  5      1    yes    r212i0n11        1
  6      2    yes    r212i0n11        2
  7      3    yes    r212i0n11        3
```

The second problem is that, as demonstrated with the ps command below, the OpenMP threads are also placed on the same core where the associated MPI process is running:

```
ps -C a.out -L -opsr,comm,time,pid,ppid,lwp

PSR COMMAND      TIME PID  PPID LWP
  0 a.out    00:00:02 4098 4092 4098
  0 a.out    00:00:02 4098 4092 4108
  0 a.out    00:00:02 4098 4092 4110
  1 a.out    00:00:03 4099 4092 4099
  1 a.out    00:00:03 4099 4092 4106
  2 a.out    00:00:03 4100 4092 4100
  2 a.out    00:00:03 4100 4092 4109
  3 a.out    00:00:03 4101 4092 4101
  3 a.out    00:00:03 4101 4092 4107
```

Sample PBS script demonstrating proper placement:

```
#PBS -l select=2:ncpus=12:mpiprocs=4:model=wes

module load mpi-hpe/mpt
module load comp-intel/2015.0.090

setenv MPI_DSM_VERBOSE
setenv OMP_NUM_THREADS 2
setenv KMP_AFFINITY disabled

cd $PBS_O_WORKDIR

#the following two lines will result in identical placement

mpiexec -np 8 omplace -nt 2 -c 0,1,3,4,6,7,9,10 -vv ./a.out
#mpiexec -np 8 omplace -nt 2 -c 0-10:bs=2+st=3 -vv  ./a.out
```

Shown in the PBS stderr file, the 4 MPI processes on each node are properly distributed on the two sockets with processes 0 and 1 on CPUs 0 and 3 (first socket) and processes 2 and 3 on CPUs 6 and 9 (second socket).

```
MPI: DSM information
MPI: using dplace
grank  lrank  pinning  node name     cpuid
  0      0     yes     r212i0n10       0
  1      1     yes     r212i0n10       3
  2      2     yes     r212i0n10       6
  3      3     yes     r212i0n10       9
  4      0     yes     r212i0n11       0
  5      1     yes     r212i0n11       3
  6      2     yes     r212i0n11       6
  7      3     yes     r212i0n11       9
```

In the PBS stout file, it shows the placement of the two OpenMP threads for each MPI process:

```
omplace: This is an HPE MPI program.
omplace: placement file /tmp/omplace.file.6454
   fork skip=0  exact cpu=0-10:3
   thread oncpu=0 cpu=1 noplace=1  exact
   thread oncpu=3 cpu=4 noplace=1  exact
   thread oncpu=6 cpu=7 noplace=1  exact
   thread oncpu=9 cpu=10 noplace=1  exact
omplace: This is an HPE MPI program.
omplace: placement file /tmp/omplace.file.22771
   fork skip=0  exact cpu=0-10:3
   thread oncpu=0 cpu=1 noplace=1  exact
   thread oncpu=3 cpu=4 noplace=1  exact
   thread oncpu=6 cpu=7 noplace=1  exact
   thread oncpu=9 cpu=10 noplace=1  exact
```

To get a better picture of how the OpenMP threads are placed, using the following ps command:

```
ps -C a.out -L -opsr,comm,time,pid,ppid,lwp

PSR COMMAND      TIME PID  PPID LWP
  0 a.out    00:00:06 4436 4435 4436
  1 a.out    00:00:03 4436 4435 4447
  1 a.out    00:00:03 4436 4435 4448
  3 a.out    00:00:06 4437 4435 4437
  4 a.out    00:00:05 4437 4435 4446
  6 a.out    00:00:06 4438 4435 4438
  7 a.out    00:00:05 4438 4435 4444
  9 a.out    00:00:06 4439 4435 4439
 10 a.out     00:00:05 4439 4435 4445
```

Pinning, the binding of a process or thread to a specific core, can improve the performance of your code by increasing the percentage of local memory accesses.

Once your code runs and produces correct results on a system, the next step is performance improvement. For a code that uses multiple cores, the placement of processes and/or threads can play a significant role in performance.

Given a set of processor cores in a PBS job, the Linux kernel usually does a reasonably good job of mapping processes/threads to physical cores, although the kernel may also migrate processes/threads. Some OpenMP runtime libraries and MPI libraries may also perform certain placements by default. In cases where the placements by the kernel or the MPI or OpenMP libraries are not optimal, you can try several methods to control the placement in order to improve performance of your code. Using the same placement from run to run also has the added benefit of reducing runtime variability.

Pay attention to maximizing data locality while minimizing latency and resource contention, and have a clear understanding of the characteristics of your own code and the machine that the code is running on.

## Characteristics of NAS Systems

NAS provides two distinctly different types of systems: Pleiades, Aitken, and Electra are cluster systems, and Endeavour is a global shared-memory system. Each type is described in this section.

### Pleiades, Aitken, and Electra

On Pleiades, Aitken, and Electra, memory on each node is accessible and shared only by the processes and threads running on that node. Pleiades is a cluster system consisting of different processor types: Sandy Bridge, Ivy Bridge, Haswell, and Broadwell. Electra is a cluster system that consists of Broadwell and Skylake nodes, and Aitken is a cluster system that consists of Cascade Lake and AMD Rome nodes.

Each node contains two sockets, with a symmetric memory system inside each socket. These nodes are considered non-uniform memory access (NUMA) systems, and memory is accessed across the two sockets through the inter-socket interconnect. So, for optimal performance, data locality should not be overlooked on these processor types.

However, compared to a global shared-memory NUMA system such as Endeavour, data locality is less of a concern on the cluster systems. Rather, minimizing latency and resource contention will be the main focus when pinning processes/threads on these systems.

For more information on Pleiades, Aitken, and Electra, see the following articles:

- Pleiades Configuration Details
- Aitken Configuration Details
- Electra Configuration Details

### Endeavour

Endeavour comprises two hosts. Each host is a NUMA system that contains 32 sockets with a total of 896 cores. A process/thread can access the local memory on its socket, remote memory across sockets within the same chassis through the Ultra Path Interconnnect, and remote memory across chassis through the HPE Superdome Flex ASICs, with varying latencies. So, data locality is critical for achieving good performance on Endeavour.

Note: When developing an application, we recommend that you initialize data in parallel so that each processor core initializes the data it is likely to access later for calculation.

For more information, see Endeavour Configuration Details.

### Methods for Process/Thread Pinning

Several pinning approaches for OpenMP, MPI and MPI+OpenMP hybrid applications are listed below. We recommend using the Intel compiler (and its runtime library) and the HPE MPT software on NAS systems, so most of the approaches pertain specifically to them. You can also use the mbind tool for multiple OpenMP libraries and MPI environments.

### OpenMP codes

- Using Intel OpenMP Thread Affinity for Pinning
- Using the dplace Tool for Pinning
- Using the omplace Tool for Pinning
- Using the mbind Tool for Pinning

### MPI codes

- Setting HPE MPT Environment Variables
- Using the omplace Tool for Pinniing

- [Using the mbind Tool for Pinning](#)

## MPI+OpenMP hybrid codes

- [Using the omplace Tool for Pinning](#)
- [Using the mbind Tool for Pinning](#)

## Checking Process/Thread Placement

Each of the approaches listed above provides some verbose capability to print out the tool's placement results. In addition, you can check the placement using the following approaches.

## Use the ps Command

ps -C *executable_name* -L -opsr,comm,time,pid,ppid,lwp

In the generated output, use the core ID under the PSR column, the process ID under the PID column, and the thread ID under the LWP column to find where the processes and/or threads are placed on the cores.

Note: The ps command provides a snapshot of the placement at that specific time. You may need to monitor the placement from time to time to make sure that the processes/threads do not migrate.

## Instrument your code to get placement information

- Call the mpi_get_processor_name function to get the name of the processor an MPI process is running on
- Call the Linux C function sched_getcpu() to get the processor number that the process or thread is running on

For more information, see [Instrumenting your Fortran Code to Check Process/Thread Placement](#).

Hyperthreading is available and enabled on the Pleiades, Aitken, and Electra compute nodes. With hyperthreading, each physical core can function as two logical processors. This means that the operating system can assign two threads per core by assigning one thread to each logical processor.

Note: Hyperthreading is currently off on Aitken Rome nodes.

The following table shows the number of physical cores and potential logical processors available for each processor type:

| Processor Model | Physical Cores (N) | Logical Processors (2N) |
|---|---|---|
| Sandy Bridge | 16 | 32 |
| Ivy Bridge | 20 | 40 |
| Haswell | 24 | 48 |
| Broadwell | 28 | 56 |
| Skylake | 40 | 80 |
| Cascade Lake | 40 | 80 |

If you use hyperthreading, you can run an MPI code using $2N$ processes per node instead of $N$ process per node—so you can use half the number of nodes for your job. Each process will be assigned to run on one logical processor; in reality, two processes are running on the same physical core.

Running two processes per core can take less than twice the wall-clock time compared to running only one process per core—if one process does not keep the functional units in the core busy, and can share the resources in the core with another process.

## Benefits and Drawbacks

Using hyperthreading can improve the overall throughput of your jobs, potentially saving standard billing unit (SBU) charges. Also, requesting half the usual number of nodes may allow your job to start running sooner—an added benefit when the systems are loaded with many jobs. However, using hyperthreading may not always result in better performance.

WARNING: Hyperthreading does not benefit all applications. Also, some applications may show improvement with some process counts but not with others, and there may be other unforeseen issues. Therefore, before using this technology in your production run, you should test your applications with and without hyperthreading. If your application runs more than two times slower with hyperthreading than without, do not use it.

## Using Hyperthreading

Hyperthreading can improve the overall throughput, as demonstrated in the following example.

**Example**

Consider the following scenario with a job that uses 40 MPI ranks on Ivy Bridge. Without hyperthreading, we would specify:

#PBS -lselect=2:ncpus=20:**mpiprocs=20**:model=ivy

and the job will use 2 nodes with 20 processes per node. Suppose that the job takes 1000 seconds when run this way. If we run the job with hyperthreading, e.g.:

#PBS -l**select=1**:ncpus=20:**mpiprocs=40**:model=ivy

then the job will use 1 node with all 40 processes running on that node. Suppose this job takes 1800 seconds to complete.

Without hyperthreading, we used 2 nodes for 1000 seconds (a total of 2000 node-seconds); with hyperthreading, we used 1 node for 1800 seconds (1800 node-seconds). Thus, under these circumstances, if you were interested in getting the best wall-clock time performance for a single job, you would use two nodes without hyperthreading. However, if you were interested in minimizing resource usage, especially with multiple jobs running simultaneously, using hyperthreading would save you 10% in SBU charges.

## Mapping of Physical Core IDs and Logical Processor IDs

Mapping between the physical core IDs and the logical processor IDs is summarized in the following table. The value of $N$ is 16, 20, 24, 28, and 40 for Sandy Bridge, Ivy Bridge, Haswell, Broadwell, and Skylake/Cascade Lake processor types, respectively.

| Physical ID | Physical Core ID | Logical Processor ID |
|---|---|---|
| 0 | 0 | $0$ ; $N$ |
| 0 | 1 | $1$ ; $N+1$ |
| ... | ... | ....... |
| 0 | $N/2 - 1$ | $N/2 - 1$; $N + N/2 - 1$ |
| 1 | $N/2$ | $N/2$ ; $N + N/2$ |
| 1 | $N/2 + 1$ | $N/2 + 1$; $N + N/2 + 1$ |
| 1 | ... | ... |

| 1 | $N$ - 1 | $N$ - 1; 2$N$ - 1 |
|---|---------|-------------------|

Note: For additional mapping details, see the configuration diagrams at the for each processor type, or run the cat /proc/cpuinfo command on the specific node type.

The dplace tool binds processes/threads to specific processor cores to improve your code performance. For an introduction to pinning at NAS, see Process/Thread Pinning Overview.

Once pinned, the processes/threads do not migrate. This can improve the performance of your code by increasing the percentage of local memory accesses.

dplace invokes a kernel module to create a job placement container consisting of all (or a subset of) the CPUs of the cpuset. In the current dplace version 2, an LD_PRELOAD library (libdplace.so) is used to intercept calls to the functions fork(), exec(), and pthread_create() to place tasks that are being created. Note that tasks created internal to glib are not intercepted by the preload library. These tasks will *not* be placed. If no placement file is being used, then the dplace process is placed in the job placement container and (by default) is bound to the first CPU of the cpuset associated with the container.

## Syntax

dplace [-e] [-c *cpu_numbers*] [-s *skip_count*] [-n *process_name*] \
      [-x *skip_mask*] [-r [l|b|t]] [-o *log_file*] [-v 1|2] \
      command [*command-args*]
dplace [-p *placement_file*] [-o *log_file*] command [mpiexec -np4 a.out]
dplace [-q] [-qq] [-qqq]

As illustrated above, dplace "execs" command (in this case, without its mpiexec arguments), which executes within this placement container and continues to be bound to the first CPU of the container. As the command forks child processes, they inherit the container and are bound to the next available CPU of the container.

If a placement file is being used, then the dplace process is not placed at the time the job placement container is created. Instead, placement occurs as processes are forked and executed.

## Options for dplace

Explanations for some of the options are provided below. For additional information, see **man dplace**.

### -e and -c *cpu_numbers*

-e determines exact placement. As processes are created, they are bound to CPUs in the exact order specified in the CPU list. CPU numbers may appear multiple times in the list.

A CPU value of "x" indicates that binding should *not* be done for that process. If the end of the list is reached, binding starts over again at the beginning of the list.

-c *cpu_numbers* specifies a list of CPUs, optionally strided CPU ranges, or a striding pattern. For example:

- -c 1
- -c 2-4 (equivalent to -c 2,3,4)
- -c 12-8 (equivalent to -c 12,11,10,9,8)
- -c 1,4-8,3
- -c 2-8:3 (equivalent to -c 2,5,8)
- -c CS
- -c BT

Note: CPU numbers are *not* physical CPU numbers. They are logical CPU numbers that are relative to the CPUs that are in the allowed set, as specified by the current cpuset.

A CPU value of "x" (or *), in the argument list for the -c option, indicates that binding should not be done for that process. The value "x" should be used only if the -e option is also used.

Note that CPU numbers start at 0.

For striding patterns, any subset of the characters (B)lade, (S)ocket, (C)ore, (T)hread may be used; their ordering specifies the nesting of the iteration. For example, SC means to iterate all the cores in a socket before moving to the next CPU socket, while CB means to pin to the first core of each blade, then the second core of every blade, and so on.

For best results, use the -e option when using stride patterns. If the -c option is not specified, all CPUs of the current cpuset are available. The command itself (which is "execed" by dplace) is the first process to be placed by the -c *cpu_numbers*.

Without the -e option, the order of numbers for the -c option is not important.

-x *skip_mask*
      Provides the ability to skip placement of processes. The *skip_mask* argument is a bitmask. If bit *N* of *skip_mask* is set, then the *N*+1th process that is forked is not placed. For example, setting the mask to 6 prevents the second and third processes from being placed. The first process (the process named by the command) will be assigned to the first CPU. The second and third processes are not placed. The fourth process is assigned to the second CPU, and so on. This option is useful for certain classes of threaded applications that spawn a few helper processes that typically do not use much CPU time.

-s *skip_count*

>Skips the first *skip_count* processes before starting to place processes onto CPUs. This option is useful if the first *skip_count* processes are "shepherd" processes used only for launching the application. If *skip_count* is not specified, a default value of 0 is used.

-q

>Lists the global count of the number of active processes that have been placed (by dplace) on each CPU in the current cpuset. Note that CPU numbers are logical CPU numbers within the cpuset, not physical CPU numbers. If specified twice, lists the current dplace jobs that are running. If specified three times, lists the current dplace jobs and the tasks that are in each job.

-o *log_file*

>Writes a trace file to *log_file* that describes the placement actions that were made for each fork, exec, etc. Each line contains a time-stamp, process id:thread number, CPU that task was executing on, taskname and placement action. Works with version 2 only.

## Examples of dplace Usage

## For OpenMP Codes

#PBS -lselect=1:ncpus=8

#With Intel compiler versions 10.1.015 and later,
#you need to set KMP_AFFINITY to disabled
#to avoid the interference between dplace and
#Intel's thread affinity interface.

setenv KMP_AFFINITY disabled

#The -x2 option provides a skip map of 010 (binary 2) to
#specify that the 2nd thread should not be bound. This is
#because under the new kernels, the master thread (first thread)
#will fork off one monitor thread (2nd thread) which does
#not need to be pinned.

#On Pleiades, if the number of threads is less than
#the number of cores, choose how you want
#to place the threads carefully. For example,
#the following placement is good on Harpertown
#but not good on other Pleiades processor types:

dplace -x2 -c 2,1,4,5 ./a.out

To check the thread placement, you can add the -o option to create a log:

dplace -x2 -c 2,1,4,5 -o *log_file* ./a.out

Or use the following command on the running host while the job is still running:

ps -C a.out -L -opsr,comm,time,pid,ppid,lwp > placement.out

## Sample Output of log_file

```
timestamp       process:thread cpu taskname| placement action
15:32:42.196786 31044       1 dplace   | exec ./openmp1, ncpu 1
15:32:42.210628 31044:0     1 a.out    | load, cpu 1
15:32:42.211785 31044:0     1 a.out    | pthread_create thread_number 1, ncpu -1
15:32:42.211850 31044:1     - a.out    | new_thread
15:32:42.212223 31044:0     1 a.out    | pthread_create thread_number 2, ncpu 2
15:32:42.212298 31044:2     2 a.out    | new_thread
15:32:42.212630 31044:0     1 a.out    | pthread_create thread_number 3, ncpu 4
15:32:42.212717 31044:3     4 a.out    | new_thread
15:32:42.213082 31044:0     1 a.out    | pthread_create thread_number 4, ncpu 5
15:32:42.213167 31044:4     5 a.out    | new_thread
15:32:54.709509 31044:0     1 a.out    | exit
```

## Sample Output of placement.out

```
PSR COMMAND        TIME  PID  PPID  LWP
 1 a.out     00:00:02 31044 31039 31044
 0 a.out     00:00:00 31044 31039 31046
 2 a.out     00:00:02 31044 31039 31047
 4 a.out     00:00:01 31044 31039 31048
 5 a.out     00:00:01 31044 31039 31049
```

Note: Intel OpenMP jobs use an extra thread that is unknown to the user and it does not need to be placed. In the above example, this extra thread (31046) is running on core number 0.

## For MPI Codes Built with HPE's MPT Library

With HPE's MPT, only 1 shepherd process is created for the entire pool of MPI processes, and the proper way of pinning using dplace is to skip the shepherd process.

Here is an example for Endeavour:

```
#PBS -l ncpus=8
....
 mpiexec -np 8 dplace -s1 -c 0-7 ./a.out
```

On Pleiades, if the number of processes in each node is less than the number of cores in that node, choose how you want to place the processes carefully. For example, the following placement works well on Harpertown nodes, but not on other Pleiades processor types:

```
#PBS -l select=2:ncpus=8:mpiprocs=4 ... mpiexec -np 8 dplace -s1 -c 2,4,1,5 ./a.out
```

To check the placement, you can set MPI_DSM_VERBOSE, which prints the placement in the PBS stderr file:

```
#PBS -l select=2:ncpus=8:mpiprocs=4
...
setenv MPI_DSM_VERBOSE
mpiexec -np 8 dplace -s1 -c 2,4,1,5 ./a.out
```

## Output in PBS stderr File

```
MPI: DSM information
grank  lrank  pinning  node name    cpuid
  0     0    yes    r75i2n13      1
  1     1    yes    r75i2n13      2
  2     2    yes    r75i2n13      4
  3     3    yes    r75i2n13      5
  4     0    yes    r87i2n6       1
  5     1    yes    r87i2n6       2
  6     2    yes    r87i2n6       4
  7     3    yes    r87i2n6       5
```

If you use the -o *log_file* flag of dplace, the CPUs where the processes/threads are placed will be printed, but the node names are not printed.

```
#PBS -l select=2:ncpus=8:mpiprocs=4
....
mpiexec -np 8 dplace -s1 -c 2,4,1,5 -o log_file ./a.out
```

## Output in log_file

```
timestamp        process:thread cpu taskname | placement action
15:16:35.848646 19807        - dplace    | exec ./new_pi, ncpu -1
15:16:35.877584 19807:0      - a.out     | load, cpu -1
15:16:35.878256 19807:0      - a.out     | fork -> pid 19810, ncpu 1
15:16:35.879496 19807:0      - a.out     | fork -> pid 19811, ncpu 2
15:16:35.880053 22665:0      - a.out     | fork -> pid 22672, ncpu 2
15:16:35.880628 19807:0      - a.out     | fork -> pid 19812, ncpu 4
15:16:35.881283 22665:0      - a.out     | fork -> pid 22673, ncpu 4
15:16:35.882536 22665:0      - a.out     | fork -> pid 22674, ncpu 5
15:16:35.881960 19807:0      - a.out     | fork -> pid 19813, ncpu 5
15:16:57.258113 19810:0      1 a.out     | exit
15:16:57.258116 19813:0      5 a.out     | exit
15:16:57.258215 19811:0      2 a.out     | exit
15:16:57.258272 19812:0      4 a.out     | exit
15:16:57.260458 22672:0      2 a.out     | exit
15:16:57.260601 22673:0      4 a.out     | exit
15:16:57.260680 22674:0      5 a.out     | exit
15:16:57.260675 22671:0      1 a.out     | exit
```

## For MPI Codes Built with MVAPICH2 Library

With MVAPICH2, 1 shepherd process is created for each MPI process. You can use ps -L -u *your_userid* on the running node to see these processes. To properly pin MPI processes using dplace, you cannot skip the shepherd processes and must use the following:

```
mpiexec -np 4 dplace -c2,4,1,5 ./a.out
```

Process/Thread Pinning

The mbind utility is a "one-stop" tool for binding processes and threads to CPUs. It can also be used to track memory usage. The utility, developed at NAS, works for for MPI, OpenMP, and hybrid applications, and is available in the /u/scicon/tools/bin directory on Pleiades.

Recommendation: Add /u/scicon/tools/bin to the PATH environment variable in your startup configuration file to avoid having to include the entire path in the command line.

One of the benefits of mbind is that it relieves you from having to learn the complexity of each individual pinning approach for associated MPI or OpenMP libraries. It provides a uniform usage model that works for multiple MPI and OpenMP environments.

Currently supported MPI and OpenMP libraries are listed below.

MPI:

- HPE-MPT
- MVAPICH2
- INTEL-MPI
- OPEN-MPI (including Mellanox HPC-X MPI)
- MPICH

Note: When using mbind with HPE-MPT, it is highly recommended that you use MPT 2.17r13, 2.21 or a later version in order to take full advantage of mbind capabilities.

OpenMP:

- Intel OpenMP runtime library
- GNU OpenMP library
- PGI runtime library
- Oracle Developer Studio thread library

Starting with version 1.7, the use of mbind is no longer limited to cases where the same set of CPU lists is used for all processor nodes. However, as in previous versions, the same number of threads must be used for all processes.

WARNING: The mbind tool might not work properly when used together with other performance tools.

## Syntax

#For OpenMP:
mbind.x [-*options*] program [*args*]

#For MPI or MPI+OpenMP hybrid which supports mpiexec:
mpiexec -np nranks mbind.x [-*options*] program [*args*]

To find information about all available options, run the command mbind.x -help.

Here are a few recommended mbind options:

-cs, -cp, -cc;
or -c*cpulist*
    -cs for spread (default), -cp for compact, -cc for cyclic; -c*cpulist* for process ranks (for example, -c*0,3,6,9*). CPU numbers in the *cpulist* are relative within a cpuset, if present.
    Note that the -cs option will distribute the processes and threads among the physical cores to minimize various resource contentions, and is usually the best choice for placement.
-t[*n*]
    Number of threads per process. The default value is given by the OMP_NUM_THREADS environment variable; this option overrides the value specified by OMP_NUM_THREADS.
-gm[*n*]
    Print memory usage information. This option is for printing memory usage of each process at the end of a run. Optional value [*n*] can be used to select one of the memory usage types: 0=default, 1=VmHWM, 2=VmRSS, 3=WRSS. Recognized symbolic values for [*n*]: "hwm", "rss", or "wrss". For default, environment variable GM_TYPE may be used to select the memory usage type:

- VmHWM - high water mark
- VmRSS - resident memory size
- WRSS - weighted memory usage (if available; else, same as VmRSS)

-gmc[*s:n*]
    Print memory usage every [*s*] seconds for [*n*] times. The -gmc option indicates continuous printing of memory usage at a default interval of 5 seconds. Use additional option [*s:n*] to control the interval length [*s*] and the number of printing times [*n*]. Environment variable GM_TIMER may also be used to set the [*s:n*] value.
-gmr[*list*]
    Print memory usage for selected ranks in the list. This option controls the subset of ranks for memory usage to print. [*list*] is a comma-separated group of numbers with possible range.
-l
    Print node information. This option prints a quick summary of node information by calling theclist.x utility.
-v[*n*]

Verbose flag; Option -v or -v1 prints the process/thread-CPU binding information. With [*n*] greater than 1, the option prints additional debugging information. [*n*] controls the level of details. Default is n=0 (OFF).

## Examples

## Print a Processor Node Summary to Help Determine Proper Process or Thread Pinning

In your PBS script, add the following to print the summary:

#PBS -lselect=...:model=bro

mbind.x -l

...

In the sample output below for a Broadwell node, look for the listing under column CPUs(SMTs). CPUs listed in the same row are located in the same socket and share the same last level cache, as shown in this configuration diagram.

```
Host Name        : r601i0n3
Processor Model  : Intel(R) Xeon(R) CPU E5-2680 v4 @ 2.40GHz
Processor Speed  : 1600 MHz (max 2401 MHz)
Level 1 Cache (D) : 32 KB
Level 1 Cache (I) : 32 KB
Level 2 Cache (U) : 256 KB
Level 3 Cache (U) : 35840 KB (shared by 14 cores)
SMP Node Memory   : 125.0 GB (122.3 GB free, 2 mem nodes)

Number of Sockets  : 2
Number of L3 Caches : 2
Number of Cores    : 28
Number of SMTs/Core : 2
Number of CPUs     : 56

Socket Cache  Cores        CPUs(SMTs)
0      0      0-6,8-14     (0,28)(1,29)(2,30)(3,31)(4,32)(5,33)(6,34)(7,35)(8,36)
                           (9,37)(10,38)(11,39)(12,40)(13,41)
1      0      0-6,8-14     (14,42)(15,43)(16,44)(17,45)(18,46)(19,47)(20,48)(21,49)
                           (22,50)(23,51)(24,52)(25,53)(26,54)(27,55)
```

## For Pure OpenMP Codes Using Intel OpenMP Library

Sample PBS script:

```
#PBS -l select=1:ncpus=28:model=bro
#PBS -l walltime=0:5:0

module load comp-intel

setenv OMP_NUM_THREADS 4
cd $PBS_O_WORKDIR

mbind.x -cs -t4 -v ./a.out
#or simply:
#mbind.x -v ./a.out
```

The four OpenMP threads are spread (with the -cs option) among four physical cores in a node (two on each socket), as shown in the application's stdout:

```
host: r635i7n14, ncpus: 56, nthreads: 4, bound to cpus: {0,1,14,15}
OMP: Warning #181: GOMP_CPU_AFFINITY: ignored because KMP_AFFINITY has been defined
```

The proper placement is further demonstrated in the output of the ps command below:

```
r635i7n14% ps -C a.out -L -opsr,comm,time,pid,ppid,lwp
PSR COMMAND        TIME   PID  PPID  LWP
 0 a.out        00:02:06 37243 36771 37243
 1 a.out        00:02:34 37243 36771 37244
14 a.out        00:01:47 37243 36771 37245
15 a.out        00:01:23 37243 36771 37246
```

Note: If you use older versions of Intel OpenMP via older versions of Intel compiler modules (comp-intel/2016.181 or earlier) during runtime, the ps output will show an extra thread that does not do any work, and therefore does not accumulate any time. Since this extra thread will not interfere with the other threads, it does not need to be placed.

## For Pure MPI Codes Using HPE MPT
WARNING: mbind.x disables MPI_DSM_DISTRIBUTE and overwrites the placement initially performed by MPT's mpiexec. The placement

output from MPI_DSM_VERBOSE (if set) most likely is incorrect and should be ignored.

Sample PBS script where the same number of MPI ranks are used in different nodes:

#PBS -l select=2:ncpus=28:mpiprocs=4:model=bro

module load comp-intel
module load mpi-hpe

#setenv MPI_DSM_VERBOSE

cd $PBS_O_WORKDIR

mpiexec -np 8 mbind.x -cs -v ./a.out
#or simply:
#mpiexec mbind.x -v ./a.out

On each of the two nodes, four MPI processes are spread among four physical cores (CPUs 0,1,14,15); two on each socket, as shown in the application's stdout:

host: r601i0n3, ncpus: 56, process-rank: 0 (r0), bound to cpu: 0
host: r601i0n3, ncpus: 56, process-rank: 1 (r1), bound to cpu: 1
host: r601i0n3, ncpus: 56, process-rank: 2 (r2), bound to cpu: 14
host: r601i0n3, ncpus: 56, process-rank: 3 (r3), bound to cpu: 15
host: r601i0n4, ncpus: 56, process-rank: 4 (r0), bound to cpu: 0
host: r601i0n4, ncpus: 56, process-rank: 5 (r1), bound to cpu: 1
host: r601i0n4, ncpus: 56, process-rank: 6 (r2), bound to cpu: 14
host: r601i0n4, ncpus: 56, process-rank: 7 (r3), bound to cpu: 15

Note: For readability in this article, the printout of the binding information from mbind.x is sorted by the process-rank. An actual printout will not be sorted.

Sample PBS script where different numbers of MPI ranks are used on different nodes:

#PBS -l select=1:ncpus=28:mpiprocs=1:model=bro+2:ncpus=28:mpiprocs=4:model=bro

module load comp-intel
module load mpi-hpe

#setenv MPI_DSM_VERBOSE

cd $PBS_O_WORKDIR

mpiexec -np 9 mbind.x -cs -v ./a.out
#Or simply:
#mpiexec mbind.x -v ./a.out

As shown in the application's stdout, only one MPI process is used on the first node and it is pinned to CPU 0 on that node. For each of the other two nodes, four MPI processes are spread among four physical cores (CPUs 0,1,14,15):

host: r601i0n3, ncpus: 56, process-rank: 0 (r0), bound to cpu: 0
host: r601i0n4, ncpus: 56, process-rank: 1 (r0), bound to cpu: 0
host: r601i0n4, ncpus: 56, process-rank: 2 (r1), bound to cpu: 1
host: r601i0n4, ncpus: 56, process-rank: 3 (r2), bound to cpu: 14
host: r601i0n4, ncpus: 56, process-rank: 4 (r3), bound to cpu: 15
host: r601i0n12, ncpus: 56, process-rank: 5 (r0), bound to cpu: 0
host: r601i0n12, ncpus: 56, process-rank: 6 (r1), bound to cpu: 1
host: r601i0n12, ncpus: 56, process-rank: 7 (r2), bound to cpu: 14
host: r601i0n12, ncpus: 56, process-rank: 8 (r3), bound to cpu: 15

## For MPI+OpenMP Hybrid Codes Using HPE MPT and Intel OpenMP

Sample PBS script:

#PBS -l select=2:ncpus=28:mpiprocs=4:model=bro

module load comp-intel
module load mpi-hpe

setenv OMP_NUM_THREADS 2
#setenv MPI_DSM_VERBOSE

cd $PBS_O_WORKDIR

mpiexec -np 8 mbind.x -cs -t2 -v ./a.out
#or simply:
#mpiexec mbind.x -v ./a.out

On each of the two nodes, the four MPI processes are spread among the physical cores. The two OpenMP threads of each MPI process run on adjacent physical cores, as shown in the application's stdout:

host: r623i5n2, ncpus: 56, process-rank: 0 (r0), nthreads: 2, bound to cpus: {0,1}

Process/Thread
Pinning

```
host: r623i5n2, ncpus: 56, process-rank: 1 (r1), nthreads: 2, bound to cpus: {2,3}
host: r623i5n2, ncpus: 56, process-rank: 2 (r2), nthreads: 2, bound to cpus: {14,15}
host: r623i5n2, ncpus: 56, process-rank: 3 (r3), nthreads: 2, bound to cpus: {16,17}
host: r623i6n9, ncpus: 56, process-rank: 4 (r0), nthreads: 2, bound to cpus: {0,1}
host: r623i6n9, ncpus: 56, process-rank: 5 (r1), nthreads: 2, bound to cpus: {2,3}
host: r623i6n9, ncpus: 56, process-rank: 6 (r2), nthreads: 2, bound to cpus: {14,15}
host: r623i6n9, ncpus: 56, process-rank: 7 (r3), nthreads: 2, bound to cpus: {16,17}
```

You can confirm this by running the following ps command line on the running nodes. Note that the HPE MPT library creates a shepherd process (shown running on PSR=18 in the output below), which does not do any work.

```
r623i5n2% ps -C a.out -L -opsr,comm,time,pid,ppid,lwp
PSR COMMAND   TIME  PID PPID  LWP
 18 a.out   00:00:00 41087 41079 41087
  0 a.out   00:00:12 41092 41087 41092
  1 a.out   00:00:12 41092 41087 41099
  2 a.out   00:00:12 41093 41087 41093
  3 a.out   00:00:12 41093 41087 41098
 14 a.out   00:00:12 41094 41087 41094
 15 a.out   00:00:12 41094 41087 41097
 16 a.out   00:00:12 41095 41087 41095
 17 a.out   00:00:12 41095 41087 41096
```

## For Pure MPI or MPI+OpenMP Hybrid Codes Using other MPI Libraries and Intel OpenMP

Usage of mbind with MPI libraries such as HPC-X or Intel-MPI should be the same as with HPE MPT. The main difference is that you must load the proper mpi modulefile, as follows:

- For HPC-X:

  module load mpi-hpcx

- For Intel-MPI:

  module use /nasa/modulefiles/testing
  module load mpi-intel

  Note that the Intel MPI library automatically pins processes to CPUs to prevent unwanted process migration. If you find that the placement done by the Intel MPI library is not optimal, you can use mbind to do the pinning instead. If you use version 4.0.2.003 or earlier, you might need to set the environment variable I_MPI_PIN to 0 in order for mbind.x to work properly.

*The mbind utility was created by NAS staff member Henry Jin.*